

CLAIMS

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
1. A method comprising:
developing a language model from a tuning set of information;
segmenting at least a subset of a received textual corpus and calculating a
perplexity value for each segment;
refining the language model with one or more segments of the received
corpus based, at least in part, on the calculated perplexity value for the one or more
segments.
 2. A method according to claim 1, wherein the tuning set of information
is application specific.
 3. A method according to claim 1, wherein the tuning set of information
is comprised of one or more application-specific documents.
 4. A method according to claim 1, wherein the tuning set of information
is a highly accurate set of textual information linguistically relevant to, but not
taken from, the received textual corpus.
 5. A method according to claim 1, further comprising a training set
comprised of at least the subset of the received textual corpus.
 6. A method according to claim 5, further comprising:
ranking the segments of the training set based, at least in part, on the
calculated perplexity value for each segment.

1 7. A method according to claim 1, wherein segmenting at least the
2 subset of the received corpus comprises:

3 clustering every N-items of the received corpus into a training unit, wherein
4 resultant training units are separated by gaps;

5 calculate the similarity within a sequence of training chunks on either side
6 of each of the gaps; and

7 select segment boundaries that maximize intra-segment similarity and inter-
8 segment disparity.

9
10 8. A method according to claim 7, wherein the resultant segment defines
11 a training chunk.

12
13 9. A method according to claim 7, wherein N is an empirically derived
14 value based, at least in part, on the size of the received corpus.

15
16 10. A method according to claim 7, wherein the calculation of the
17 similarity within a sequence of training units defines a cohesion score.

18
19 11. A method according to claim 10, wherein intra-segment similarity is
20 measured by the cohesion score.

21
22 12. A method according to claim 7, wherein inter-segment disparity is
23 approximated from the cohesion score.
24
25

1 **13.** A method according to claim 7, wherein the calculation of inter-
2 segment disparity defines a depth score.

3
4 **14.** A method according to claim 1, wherein the perplexity value is a
5 measure of the predictive power of a certain language model to a segment of the
6 received corpus.

7
8 **15.** A method according to claim 1, further comprising:
9 ranking the segments of at least the subset of the received corpus based, at
10 least in part, on the calculated perplexity value of each segment; and
11 updating the tuning set of information with one or more of the segments
12 from at least the subset of the received corpus.

13
14 **16.** A method according to claim 15, wherein one or more of the
15 segments with the lowest perplexity value from at least the subset of the received
16 corpus are added to the tuning set.

17
18 **17.** A method according to claim 1, further comprising:
19 utilizing the refined language model in an application to predict a likelihood
20 of another corpus.

21
22 **18.** A storage medium comprising a plurality of executable instructions
23 including at least a subset of which, when executed, implement a method
24 according to claim 1.
25

1 **19.** A system comprising:
2 a storage medium having stored therein a plurality of executable
3 instructions; and
4 an execution unit, coupled to the storage medium, to execute at least a
5 subset of the plurality of executable instructions to implement a method according
6 to claim 1.

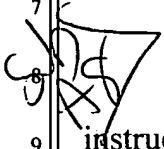
7
8 **20.** A storage medium comprising a plurality of executable instructions
9 which, when executed, implement a language modeling agent to develop a
10 language model from a tuning set of information, to segment at least a subset of a
11 received textual corpus and calculate a perplexity value for each segment, and to
12 refine the language model with one or more segments of the received corpus
13 based, at least in part, on the calculated perplexity value for the one or more
14 segments.

15
16 **21.** A storage medium according to claim 20, wherein the language
17 modeling agent utilizes a tuning set of information relevant to that of the received
18 corpus.

19
20 **22.** A storage medium according to claim 20, wherein the language
21 modeling agent ranks the segments of the training set based, at least in part, on a
22 measure of similarity between two or more segments.

23
24 **23.** A storage medium according to claim 22, wherein the similarity
25 measure is calculated for adjacent segments.

1 24. A storage medium according to claim 20, wherein the language
2 modeling agent segments the received corpus by clustering every N items of the
3 received corpus into a training unit, wherein the training units are separated by
4 gaps, calculating the similarity within a sequence of training units on either side of
5 each of the gaps, and selecting segment boundaries that improve intra-segment
6 similarity and inter-segment disparity.

7
8  25. A storage medium according to claim 20, further comprising
9 instructions to implement an application which selectively invokes the language
10 modeling agent to predict a likelihood of another corpus.

11
12 26. A storage medium according to claim 25, wherein the application is
13 one or more of a spelling and/or grammar checker, a word-processor, a speech
14 recognition application, a language translation application, and the like.

15
16 27. A system comprising:
17 a storage medium drive, to removably receive a storage medium according
18 to claim 20; and
19 an execution unit, coupled to the storage medium drive, to execute at least a
20 subset of the plurality of instructions and implement the language modeling agent.

21
22 28. A modeling agent comprising:
23 a controller, to receive invocation requests to develop a language model
24 from a corpus; and
25 a data structure generator, responsive to the controller, to develop a
language model from a tuning set of information, segment at least a subset of the

received corpus, calculate a perplexity value for each segment, and refine the language mode with one or more segments of the received corpus based, at least in part, on the calculated perplexity value.

29. A modeling agent according to claim 28, wherein the tuning set is dynamically selected as relevant to the received corpus.

30. A modeling agent according to claim 28, the data structure generator comprising:

a dynamic lexicon generation function, to develop an initial lexicon from the tuning set, and to update the lexicon with select segments from the received corpus.

31. A modeling agent according to claim 28, the data structure generator comprising:

a frequency analysis function, to determine a frequency of occurrence of segments within the received corpus.

32. A modeling agent according to claim 28, the data structure generator comprising:

a dynamic segmentation function, to iteratively segment the received corpus to improve a predictive performance attribute of the modeling agent.

1 **33.** A modeling agent according to claim 32, wherein the dynamic
2 segmentation function iteratively re-segments the received corpus until the
3 language model reaches an acceptable threshold.

4
5 **34.** A modeling agent according to claim 32, the data structure generator
6 further comprising:

7 a frequency analysis function, to determine a frequency of occurrence of
8 segments within the received corpus.

9
10 **35.** A modeling agent according to claim 34, wherein the data structure
11 generator selectively removes segments from the data structure that do not meet a
12 minimum frequency threshold, and dynamically re-segments the received corpus to
13 improve predictive capability while reducing the size of the data structure.